

Západočeská univerzita v Plzni
Fakulta aplikovaných věd
Katedra informatiky a výpočetní techniky

Diplomová práce

Nástroj pro automatickou identifikaci KIR alel

Místo této strany bude
zadání práce.

Prohlášení

Prohlašuji, že jsem diplomovou práci vypracovala samostatně a výhradně s použitím citovaných pramenů.

V Plzni dne 16. dubna 2020

Kateřina Kratochvílová

Poděkování

Ráda bych poděkovala Ing. Lucii Houdové, Ph.D. za cenné rady, věcné připomínky, trpělivost a ochotu, kterou mi v průběhu zpracování této práce věnovala. Dále bych chtěla poděkovat panu ing. Jiřímu Fatkovi za jeho rady a pomoc při vytváření praktické části.

Abstract

The text of the abstract (in English). It contains the English translation of the thesis title and a short description of the thesis.

Abstrakt

Text abstraktu (česky). Obsahuje krátkou anotaci (cca 10 řádek) v češtině. Budete ji potřebovat i při vyplňování údajů o bakalářské práci ve STAGu. Český i anglický abstrakt by měly být na stejné stránce a měly by si obsahem co možná nejvíce odpovídat (samozřejmě není možný doslovný překlad!).

Obsah

1	Úvod	8
2	Imunitní systém a jeho spojitost s geny	9
2.1	Geny	9
2.2	Imunitní systém	9
2.3	HLA a non-HLA geny	10
2.3.1	Alela a gen	11
2.4	Natural killer a jeho receptory	12
2.4.1	Natural killer	12
2.4.2	NKG2D receptor	13
2.4.3	KIR receptor	14
2.5	Nalezení vhodného dárce	20
2.6	Bordel haplotypy	21
3	Sekvenační metody získávání DNA dat	22
3.1	Sanger sequencing	22
3.2	NGS next-generation sekvenování	23
3.2.1	454 sekvenování a Ion Torrent	24
3.2.2	Illumina	25
3.2.3	SOLiD	25
3.3	Metody třetí generace	25
3.4	Read	26
3.5	Single-end, paired-end a mate-pair	26
3.6	Bordel	27
4	Analyza dostupných bioinformatických nástrojů pro zpracování NGS dat	29
4.0.1	Vytvoření testovacího haplotypu	29
4.1	ART	29
4.2	Bowtie	31
4.2.1	Burrows-Wheeler transformace	32
4.2.2	bordel	35
4.2.3	bordel	35

5	Bordel	37
5.1	ART	37
5.1.1	pokus to nějak spustit	37
5.1.2	FASTQ	37
5.1.3	bordel	39
5.2	IGV	39
6	Implementace	40
6.1	Popis problému	40
6.2	Návrh systému	40
6.3	Referenční geny	40
6.4	Použité programové prostředky	41
6.4.1	Python	41
7	Vyhodnocení výsledků a jejich srovnání	42
8	Závěr	43
9	Seznam zkratk	44
10	Výkladový slovník pojmů	45
	Literatura	46
A	Uživatelská dokumentace	49
A.1	Nastavení ART a jeho spuštění	49

1 Úvod

Transplantace krvetvorných buněk se využívá jako terapeutická procedura pro mnoho vážných hematologických poruch mezi které patří například akutní myeloidní leukemie. Transplantace jako taková je poměrně jednoduchý proces, kdy jsou dárci odebrány krvetvorné buňky a vpraveny do těla pacienta trpícím hematologickou poruchou. Problém nastává při reakci imunitního systému na nově vložený štěp. V případě, že si štěp s imunitním systémem nebudou rozumět, může dojít k silné zánětlivé reakci, která může skončit až smrtí pacienta.

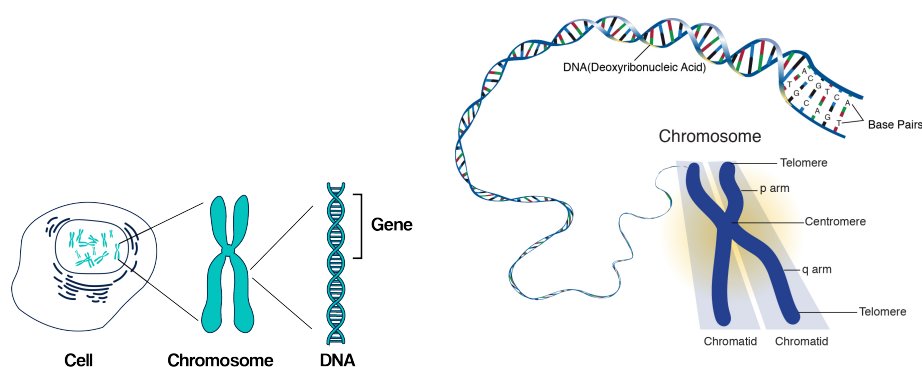
K potlačení odmítnutí se vybírají dárci podle shody v HLA znacích, věku a pohlaví. Ovšem ani to není bezrizikové. V poslední době se množí studie, které prokazují vliv takzvaných non-HLA genů. Jedním z nich může být i gen Killer-cell immunoglobulin-like receptor (KIR). V případě, kdy by se rozhodovalo mezi více dárce by se mohl ten vhodnější vybrat právě na základě KIR. Pro zjištění jak HLA znaků tak KIR genů se využívají sekvenční metody. [20]

Cílem práce je navrhnout a implementovat nástroj pro automatickou identifikaci KIR alel. Vstupní data tzv. ready jsou neznámý kus DNA (posloupnost písmen A, C, G a T) a jsou výstupem ze sekvenčních metod. Tyto data budou pro vývoj nástroje simulována nástrojem ART a v konečné fázi testování budou data vyměněna za data z FN Plzeň. Jelikož je třeba odhadnout co se pod danou posloupností nachází, bude použit nástroj bowtie2 pro zarování readů vzhledem k referenčním KIR genům. V poslední fázi bude vyhodnocena shoda readů a referenčních genů.

2 Imunitní systém a jeho spojitost s geny

2.1 Geny

V každé buňce lidského organismu, konkrétně v buněčném jádře, je možné nálezt 46 chromozomů. Jeden chromozom představuje stočenou dlouhou molekulu DNA (Deoxyribonukleovou kyselinu). Všechny 46 chromozomů obsahuje okolo 100 000 genů. Drobný segment DNA, který řídí buněčnou funkci je právě gen. Konkrétní forma genu je alela. [26]



Obrázek 2.1: Převzato z [4] a [1]

Uvnitř buňky máme celý genom který se ovšem nemusí projevit na povrchu buňky. Pokud se vlastnost kterou gen přenáší projeví na povrchu buňky označujeme to jako exprese genu (jeho sebevyjádření). Od toho se odvíjí i konkrétní názvosloví typu KIR gen, KIR receptor či molekula.

2.2 Imunitní systém

Imunitní systém chrání organismus před škodlivinami. Skládá se ze dvou hlavních částí vrozené imunity a získané imunity. Reakce imunitního systému je vždy komplexní reakce organismu mezi jednotlivými buňkami imunitního systému reagující na přítomnost specifických antigenů. Antigeny jsou látky, které imunitní systém rozpozná a zareaguje na ně. V podstatě to může být jakákoli bílkovina sloučenina. Antigen se obvykle nachází na povrchu buňky jako vyjádření genu. Imunitní systém následně zjistí o jaký antigen se jedná,

respektivě o jakou buňku se jedná, zda tělu vlastní (např. zdravá buňka) nebo buňku tělu cizí (např. nádorová buňka), tedy jedná-li se o exprese lidského genu nebo například viru. Jedná-li se o buňku tělu cizí imunitní systém reaguje snahou ji zničit.

Vrozená imunita též označována přirozená, neadaptivní, antigenně nespecifická je neměnně zapsána v DNA. To znamená, že při každém setkání s antigenem odpoví stejnou reakcí. Buňky nesoucí vrozenou imunitu jsou stále přítomně v krvi, takže jejich případná aktivace je takřka okamžitá (minuty až hodiny). Do této imunity patří i natural killer buňky s KIR receptory, které budou dále rozebírány v textu.

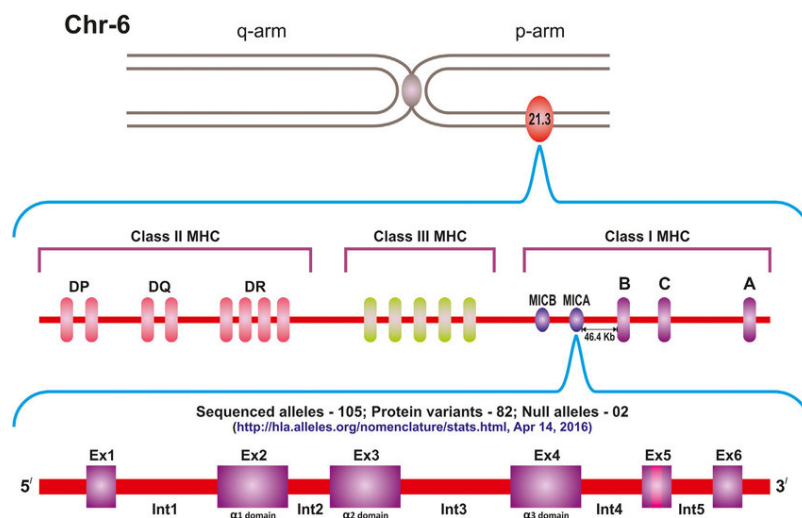
Získaná imunita též označována specifická či adaptivní oproti specifické má v genomu zapsány pouze své základy. V průběhu lidského života se vyvíjí a mění. Změna může nastat například očkování nebo proděláním patřičné choroby. Tato změna ovšem nemusí být trvalá. Z těchto důvodů může být odpověď získané imunity při setkání se stejnou chorobou rozdílná. Fungování získané imunity zajišťují T- a B- lymfocyty, ale nefunguje samostatně. Při zabíjení patogenů spolupracuje s vrozenou imunitou.

2.3 HLA a non-HLA geny

Human leucocyte antigen (HLA) je genetický systém, který je primárně zodpovědný za rozeznávání vlastního od cizorodého. Někdy je termín HLA zaměňován s MHC. MHC (Major histocompatibility complex) je souhrnný termín pro všechny komplexy, kdy podskupinou jsou právě HLA (H - Human) který je pro lidi. Stejně tak existuje DLA (D - Dog) který je pro psy. Z funkčního i biologického hlediska jde však u všech savců o stejnou skupinu genů. [20]

Přesná definice mezi HLA a non-HLA geny neexistuje. Mimo jiné i jejich rozdělení není v literaturách sjednocené. Jak je vidět z obrázku 2.2 je možné geny rozdělit do tří tříd. V některých literaturách je možné nalést označení non-HLA genů jako geny III. třídy v jiné, že jsou to všechny geny III třídy a některé geny třídy I. Tato práce se bude v označení za gen non-HLA či HLA odkazovat na hla.alleles [24]. Zjednodušeně tedy můžeme říci, že geny které nejsou řazeny k HLA skupinám jsou non-HLA. Je-li gen označen za non-HLA neznamena to, že by neměl souvislost s funkcí imunitního systému. Naopak má, jen ne výlučně s HLA systémem. Non-HLA geny kódují pro-

dukty spojené s imunitními procesy. Mezi non-HLA geny mimo jiné patří MICA, MICB a KIR. [24]



Obrázek 2.2: Šestý chromozom zobrazující HLA i non-HLA geny. Protein vzniklý expresí MICA genu je definován exony, které definují přepis do RNA. Introny v praxi nehrají roli a často jsou sekvenovány jen exony. [7]

HLA a některé non-HLA geny se nacházejí na krátkém raménku 6 chromozomu, konkrétně 6p21.3 a zaujímá úsek přibližně jednu tisícinu genomu. Tento region je nejvíce komplexní a polymorfní na lidském genomu s více než 220 geny. Oproti tomu jedna ze skupin non-HLA genů, konkrétně KIR geny, se nachází na 19 chromozomu. Rozsáhlá diverzita genů vznikala snahou eliminovat neustále se měnící spektrum patogenů. Produkty těchto genů na povrch buňky významně ovlivňují odpověď na infekční choroby a výsledky buněčné či orgánové transplantace. [24]

2.3.1 Alela a gen

Alelu můžeme definovat jako variantu genu s nepatrným rozdílem v sekvenci nukleotidů DNA oproti jiné alele stejného genu. Geny se vyskytují minimálně ve dvou formách (dvou alelách), mnohdy jich, ale může být více. U jednoho člověka mohou být přítomny pouze dvě rozdílné alely daného genu. Gen určuje výskyt nějaké vlastnosti, například tento živočich bude mít oči. Alela pak určuje jakou barvu budou mít.

V případě genu KIR2DL1 mohou být jeho alely 0010101 a 0010102. Zápis genů tak, jak s nimi budeme pracovat může vypadat způsobem zobrazeným

v 2.3.1.

$$\begin{aligned} > KIR : KIR00001 KIR2DL1 * 0010101 14738 bp \\ GTTCGGGAGGTTGGATCTCAGACGTG... \end{aligned} \quad (2.3.1)$$

TODO: Když najdu novou sekvenci tak kde je rozdíl jestli je to nový gen nebo nová alela? Není to tak že na daný pozici v genomu je vždycky gen.. a alela určuje tu vlastnost? A na co je mi teda lotus? geny jsou již plně definované - The Human Genome Project (HGP) <https://www.genome.gov/human-genome-project/What> (geny jsou ty 2DL1, 3DL1....)

jde o nové varianty - alelické skupiny, konkrétní alely - to je ve vazbě na to, jaký protein je kódován

TODO tohle je asi jen HLA nevím jestli existuje něco jako obecné rozdělení genu a aleli, možná že rozdíl bude jen v tom že těch čísel pak může být za hvězdičkou několik v závislosti o alelu jaké skupiny genů se jedná Aleli jdou definovány HLA-DRB1* což označuje označuje lokus, následované 4 čísly. TODO nevím jestli mám nějak rozebírat to že je tam HLA-DRB1 že tam je tam jednička na konci, já totiž nevím co to znamená

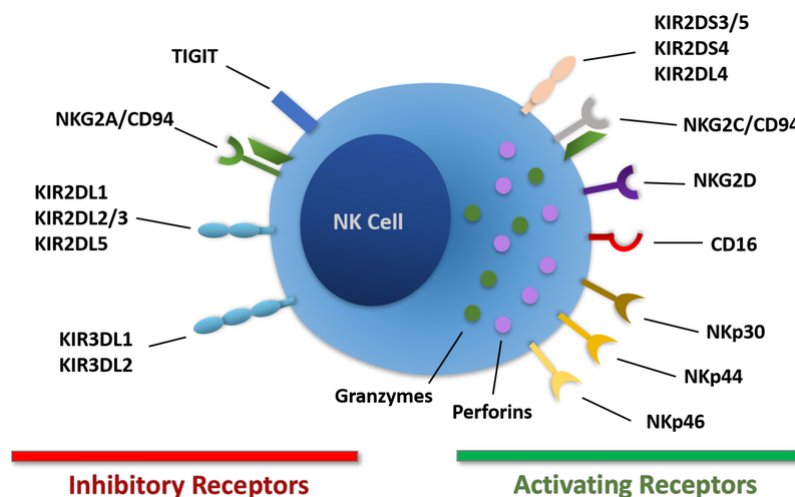
TODO tohle nevím jestli tam dávat: Alela zajišťuje konkrétní fenotypový projev genu. U jedince mohou na homologních jaderných chromozomech být přítomny pouze dvě alely. Když jsou v párových lokusech obě alely shodné, jde buď o dominantního homozygota (AA) nebo o recesivního homozygota (aa). Když jsou na párových chromozomech v daném lokusu přítomny různé alely, jde o heterozygota (Aa). Značení alel vzniká dohodou.

2.4 Natural killer a jeho receptory

2.4.1 Natural killer

Natural killer buňky (NK buňky) jsou velké granulární lymfocyty vrozeného imunitního systému. V krevním oběhu lidského těla je jich možné nalést 10–15%. Klíčovou vlastností NK buněk je nejenom schopnost rozlišit poškozené buňky od zdravích, ale i poškozené buňky rychle a efektivně likvidovat. Poškozené buňky mohou být buňky infokované virem či buňky transformované v nádorové. Na povrchu NK buňky se nachází receptory, které jsou zobrazeny na obrázku 2.3, regulující odpověď imunitního systému. Natural killer buňky oproti B- a T- lymfocitům (buňkám získané imunity) nemají

antigenně specifické receptory. Jedním ze způsobů jak NK buňky rozpoznávají a zabíjejí poškozené buňky je na základě interakce mezi KIR receptorem a HLA molekulou na povrchu zkoumané buňky (podrobněji viz sekce KIR). Stejně tak mohou zabíjet na základě receptoru NKG2D, který aktivuje cytotoxickou reakci při setkání s ligandem MICA a MICB. Ligandem označujeme malou molekulu, která se váže na vazebné místo cílového proteinu(receptoru) a vyvolává fyziologickou odpověď která může mít inibiční či aktivační charakter.



Obrázek 2.3: Natural killer buňka a její receptory, rozděleny na aktivační a inibiční. Pro tuto práci jsou důležité hlavně KIR receptory a NKG2D. [8]

2.4.2 NKG2D receptor

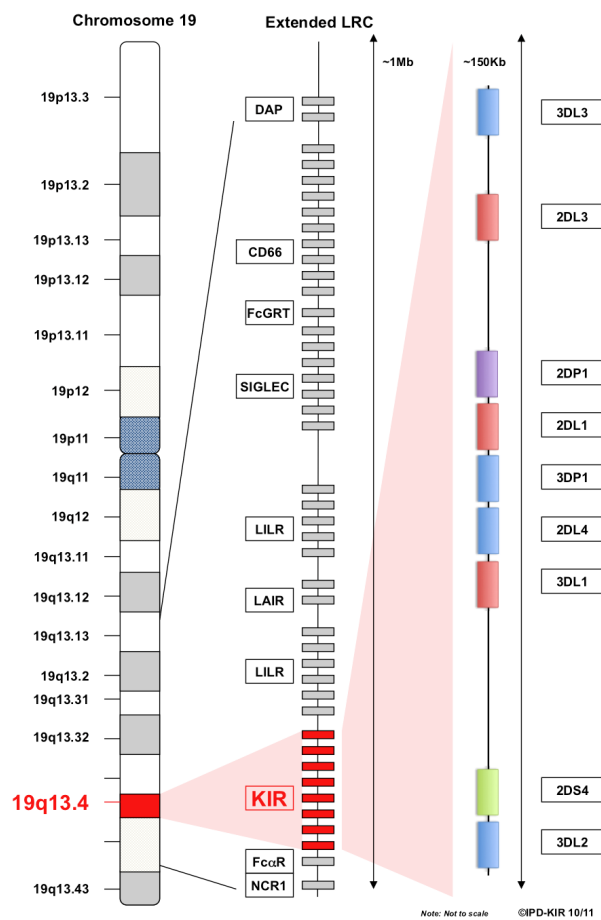
NKG2D je jeden z nejvýznamnějších aktivačních receptorů na NK buňce rozpoznávající především buněčný stres, který může spustit cytotoxicitu (schopnost ničit buňky) i když se na povrchu buňky nachází inibiční HLA-I ligandy.

Geny skupiny MICA a MICB jsou označeny jako class I chain-related gene. To znamená, že se běžně neřadí do I. třídy MHC. Takto označované geny mají souvislost s MHC I třídy, ale narozdíl od nich neváží peptidy. Oproti HLA genům, které mají svoje produkty na lymfocytech, se produkty MICA a MICB nachází na epitelových buňkách. Nejedná se tedy o standardní HLA geny, proto jsou nověji v literaturách označovány jako non-HLA. Jejich expresí na povrch buňky jsou ligandy, které se váží na receptor NKG2D. Buňky s ligandy MICA a MICB se množí při nádorovém onemocnění, zanětu nebo pod vlivem různých forem buněčného stresu a díky navázáním na receptor

může být spuštěna imunitní reakce. [22] [11] [8] [24]

2.4.3 KIR receptor

Killer immunoglobulin-like receptor (KIR) je skupina genů řazených mezi non-HLA geny. Jejich zvláštností je fakt, že se nenachází na 6 chromozomu, ale na 19 a tak shodní dárci HLA znaků mohou být neshodní v KIR znacích. Jejich expresí jsou receptory na povrchu natural killer buněk. Dnes je známo 15 genů a 2 pseudogeny rozlišujících se na inhibiční a aktivační na základě cytoplasmatického ocásku a počtu imunoglobulinových domén. [20]



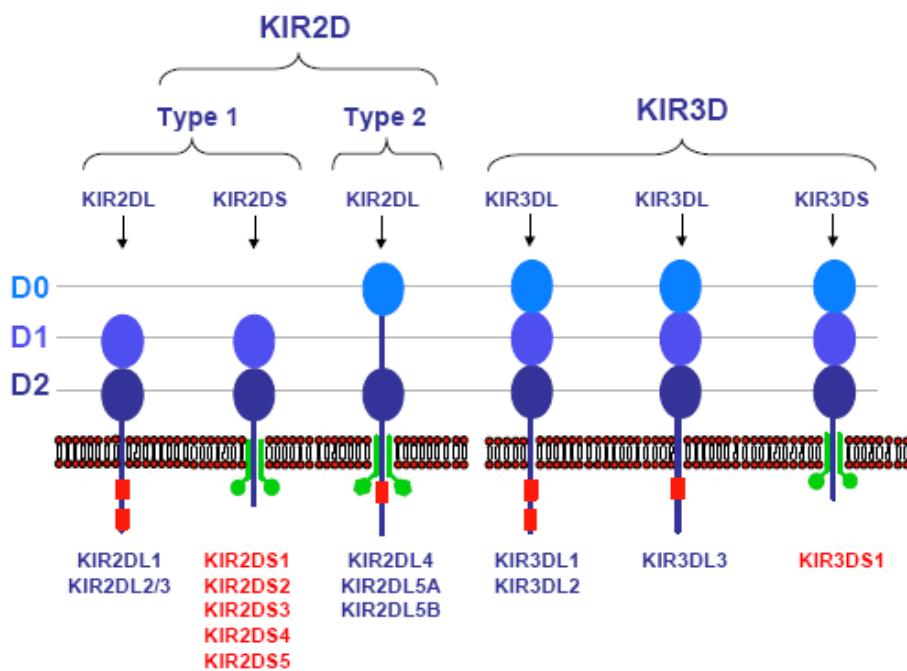
Obrázek 2.4: KIR se nachází na 19 chromozomu v oblasti jménem leukocyte receptor complex (LRC). [24]

Nomenklatura KIR genů

KIR geny (na obrázku 2.5) se liší různou délkou cytoplasmatických ocásku (tail) a různým počtem imunoglobulin-like domén (lg-like). Na základě této

rozmanitosti byla založena nomenklatura KIR genů, tedy jejich pojmenování.

Jak je vidět na obrázku 2.5, cytoplasmatický ocásek může být dlouhý (long - L) nebo krátký (short - S). Oproti tomu imunoglobulinové domény se mohou vyskytovat 2 (2D) nebo 3 (3D). Právě z těchto vlastností vychází základ pojmenování KIR genů. Příkladem může být KIR2DL1*010101, kde 2D označuje dvě imunoglobulinové domény, L značí dlouhý ocásek, 1 značí že je to první 2DL protein. Následuje hvězdička oddělující gen od alely. První tři čísla označují alely, které se liší v sekvencích jejich kódovaných proteinů, další dvě číslice se používají k rozlišení alel, které se liší synonymními rozdíly v kódující sekvenci. Konečné dvě cifry rozlišují alely na základě substituce v intronu, promotoru nebo jiné nekódující oblasti. [24]

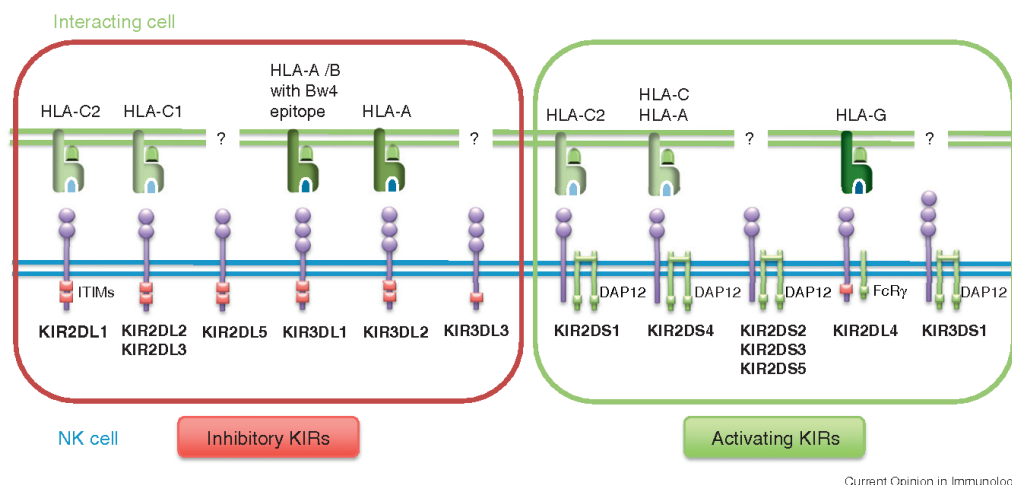


Obrázek 2.5: Nomenklatura KIR genů. [20]

Další rozdělení KIR genů je již výše zmíněné inhibiční a aktivační. Na obrázku 2.5 je možné si povšimnout detailu, že až na KIR2DL4 jsou aktivační KIR s krátkým ocáskem, zatímco inhibiční jsou s dlouhým ocáskem.

Aktivace NK buněk pomocí KIR

Jak již bylo výše zmíněno KIR receptory můžeme rozdělit na inhibiční a aktivační. Zda dojde k aktivaci NK buňky rozhoduje právě jejich rovnováha na zkoumané buňce. Zatímco inhibiční receptory se váží hlavně na molekuly HLA, aktivační receptory rozpoznávají molekuly které jsou exprimovány na membránu při buněčném stresu. Obrázek 2.6 uvádí vazebné ligandy pro jednotlivé KIR receptory.

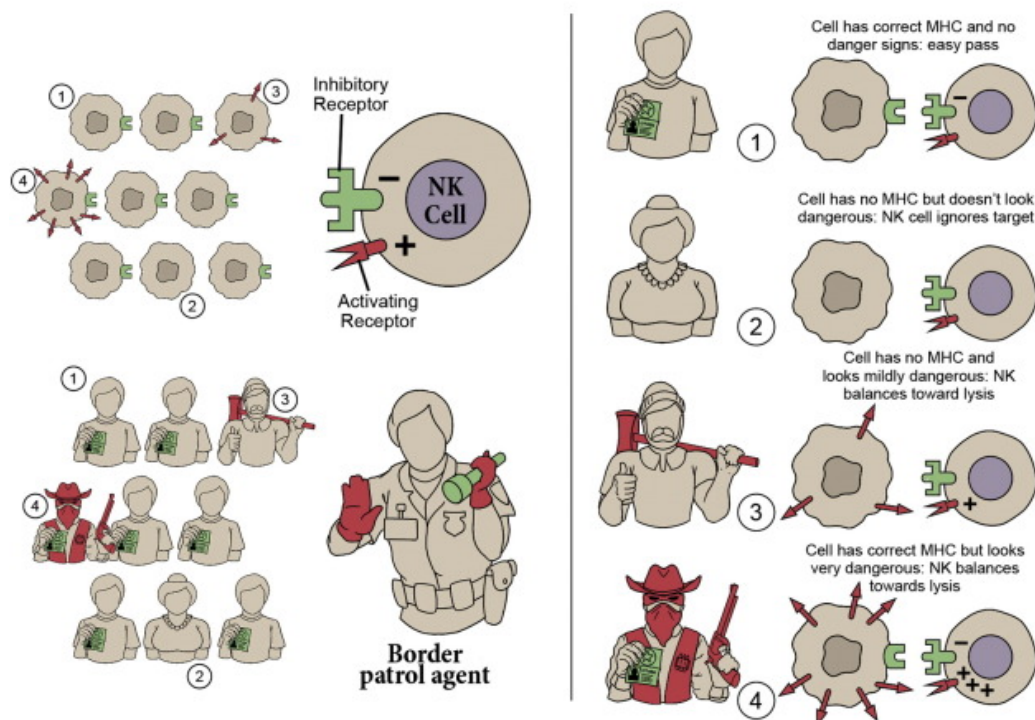


Obrázek 2.6: KIR geny a jejich vazebné ligandy. Pokud je v obrázku ? značí to, že pro daný receptor není znám vazebný ligand. [27]

NK buňky ustavičně prohledávají své okolí a testují přítomnost příslušných HLA ligand pro své KIR receptory. Pokud je příslušný HLA ligand přítomen naváže se na NK buňku (2.7 případ 1). Tímto systémem jsou ochráněny vlastní buňky. Pokud přítomen není je spuštěna cytotoxická reakce a zkoumaná buňka je zničena.

Některé virem napadené buňky potlačují propsání HLA ligand na povrch buňky a tím se brání cytotoxicitě proti T lymfocitům, ale naopak jsou více citlivější na cytotoxicitu proti NK buňkám, jak je zobrazeno na obrázku 2.7 případ 3.

The NK Cell is like a border patrol agent



Obrázek 2.7: Přirovnání fungování natural killer buňky k pasové kontrole. V pravé části jsou zobrazené případy které mohou nastat když natural killer buňka potká jinou buňku. V prvním případě je tělu vlastní zdravá buňka, kde se KIR receptor naváže na HLA ligand a k cytotoxické reakci nedojde. Druhým případem je červená krvinka. K reakci NK buňky opět nedojde, protože na zkoumané buňce nepřevažují aktivační receptory. V 3 případě je to nádorová buňka, která schová HLA ligand (může nastat po transplantaci kostní dřeně) a tím se "schová" proti T- lymfocytům. Avšak aktivační receptory převažují a tak k cytotoxicitě dojde. Ve 4 příkladě je nádorová buňka nebo virem nakažená buňka (stresové ligandy). Aktivační receptory převažují k cytotoxicitě dojde.[25]

KIR haplotyp

KIR haplotyp je vyjádření jaké konkrétní KIR geny genom obsahuje. Doposud nebylo zavedeno konkrétní pravidlo na jejich pojmenovávání. Avšak bylo navrženo, aby každý KIR haplotyp byl označen "KH – " následovaným trojmístným číslem, které bude označovat konkrétní haplotyp. Bylo by tak možné pojmenovat 999 haplotypů. [24]

Dále by se haplotypy rozdělovali na dvě skupiny A a B. Skupina B musí

obsahovat alespoň jeden z genů KIR2DL5, KIR2DS1, KIR2DS2, KIR2DS3, KIR2DS5 a KIR3DS1. Naopak skupina A neobsahuje ani jeden z těchto genů. Z tohoto pravidla je patrné, že haplotypy B mají vždy více aktivačních KIR než haplotypy A. Za trojmístným číslem by tedy dále bylo písmeno A nebo B.

Nakonec by byl připojen 17-ti místný binární kód, který by označoval přítomnost "1" či absenci "0" genu. Pořadí genů by odpovídalo pořadí v genomu od centrometrické části k telemetrické části.

Výsledné pojmenování by mohlo vypadat následovně:

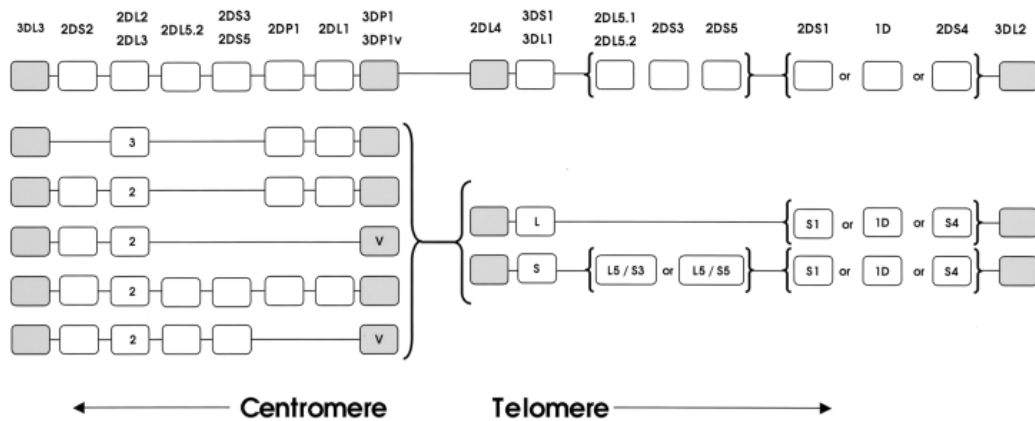
$$KH - 001A - 11100010011011011 \quad (2.4.1)$$

Je třeba si zde uvědomit, že každý jedinec má 2 KIR haplotypy. Je tedy možné dostat 4 kombinace - A/A, A/B, B/A nebo B/B. Haplotyp jedince je označován za A v případě kdy má kombinaci A/A a za B v případě jedné z kombinací A/B, B/A nebo B/B. Je možné si povšimnout, že u haplotypu B převládají inhibiční KIR geny a proto jsou dárci lépe přijímáni.

Hapl Group	Genotype ID ¹	3DL1	2DL1	2DL3	2DS4	2DL2	2DL5	3DS1	2DS1	2DS2	2DS3	2DS5	2DL4	3DL2	3DL3	2DP1	3DP1	Populations	Individuals
AA	1																	190	7,540
Bx	2																	178	2,522
Bx	4																	178	2,096
Bx	3																	167	1,157
Bx	5																	161	1,536
Bx	6																	155	899
Bx	7																	134	583
Bx	8																	130	635
Bx	9																	120	395
Bx	71																	112	443

Obrázek 2.8: Deset nejčastější KIR haplotypů. Šedý obdelník značí přítomnost genu, bílý jeho nepřítomnost. [6]

Na základě variací obsahu genů by bylo možné vytvořit nepřeborné množství KIR genotypů. Na základě sesbíraných haplotypů byl sestaven model, který toto množství mírně redukuje. Haplotyp se rozděluje na dvě části na centrometickou a telemetrickou. Kdy jednotlivé části mezi sebou mohou být kombinovány. Existují vzácné varianty, které se do tohoto modelu nehodí. [12]



Obrázek 2.9: Rozdělení KIR genů na centrometrickou a telometrickou část, pojmenování je na základě, zda je úsek blíže k centromeru nebo k telomeru (viz obrázek 2.1). [12]

Centrometická polovina je charakterizována přítomností jednoho z 2DL3 nebo 2DL2, vzácně nemusí být přítomný ani jeden. V případě 2DL2 je následně přítomen 2DS2. Tento pár genů se následně objevu v kombinaci s -2DP1, -2DL1 a -3DP1. 2DL5 gen je v centromerické části párován s 2DL2 a 2DS3, ve vzácných případech se může objevit i s 2DL2 a 2DS5. Oproti tomu při přítomnosti KIR2DL3 se dále vyskytuje KIR2DP1, -2DL1 a -3DLP1.

Telometrická polovina haplotypu je charakterizována přítomností jednoho z 3DL1 nebo 3DS1, vzácně nemusí být přítomný ani jeden. Gen 3DL1 se následně objevuje v přítomnosti s 2DS4, 1D nebo 3DL2. V případě KIR3DL se jedná o takzvaný krátký segment obsahující 2DS4 nebo KIR1D zakončené 3DL2. V případě KIR3DS1 se jedná o dlouhý segment obsahující 2DL5, párováný s 2DS3 nebo 2DS5, následovaný 2DS1, 2DS4 nebo KIR1D opět zakončený 3DL2. [12]

Podle některých studií zabývajících se vlivem KIR haplotypů na výsledky transplantace bylo zjištěno, že KIR haplotypy ovlivňují výsledky u akutní myeloidní leukémie. Ve srovnání s haplotypem A měl haplotyp B, především jeho centrometická část, ochranný účinek před návratem nemoci a zároveň zvýšil pravděpodobnost přežití pacienta. Na základě této skutečnosti se mohou dárce řadit do tří skupin best, better a neutral. Best je definován jako Cen-B/B a Tel-x/x, better jako Cen-A/x a Tel-B/x, a netral v případě jedné B části nebo žádné. [9]

Dále je možné se setka s takzvaným B-skóre

TODO podle tohoto se právě určuje B-content score pro definování

úrovně B haplotypu, best varianta pro transplantaci je B/B (cen) s B/B (tel), viz další komentář

We defined the KIR B-content score for each donor's KIR genotype as the number of centromeric and telomeric gene-content motifs containing B haplotype-defining genes. Permissible values for the KIR B-content score are 0, 1, 2, 3, and 4 (Figure 1C). A calculator for classification of the donor KIR B status (best, better, neutral) may be found at <http://www.ebi.ac.uk/ipd/kir/>.

TODO to znamená že to b score je to best, better, neutral a ještě tam je přidáný none of this?

2.5 Nalezení vhodného dárce

Mezi rizika při transplantaci krvetvorných buněk patří reakce štěpu proti hostiteli nebo relaps onemocnění (návrat nemoci). Ač je dárce vybírán podle shody v HLA znacích, sekundární kritéria jako jsou pohlaví a věk hrají také roli pro úspěšnost transplantace. Navíc podle nedávných studií výsledky příjetí štěpu ovlivňují nejenom HLA geny ale i non-HLA geny. Jedním z nich může být právě killer immunoglobulin-like receptor (KIR). V případě kdy by bylo nalezeno více vhodných dárců, tj. se shodou 10/10 nebo 9/10, vybíralo by se následně podle KIR genů. [20] [10]

Při určování shody dárce a pacienta se rozhoduje na základě shody alel u genů HLA -A, -B, -C, -DRB1, -DQB1. Díky velké diverzitě HLA genů je počet možných kombinací několik miliard. Některé kombinace genů se vyskytují na základě oblasti či národnosti častěji nebo mohou být naopak vzácné. HLA geny se obvykle dědí jako blok (celý haplotyp), avšak ve výjimečných případech může dojít k rekombinaci. Z tohoto důvodu je nejsnadnější nalést shodu v pokrevním příbuzenstvu.

Jelikož každý jedinec má dvakrát geny na pozicích HLA -A, -B, -C, -DRB1 a -DQB1 (jednu pětici od otce, druhou pětici od matky), je maximální shoda 10/10 (shoda obou alel v lokusech). Čím je shoda menší tím větší je riziko nepřijetí štěpu. U nepříbuzných jedinců lze tolerovat shodu 9/10 či 8/10. [10] [20]

V posledních letech se objevuje Haploidentická transplantace, kdy je možné použít krvetvorné buňky příbuzného se shodou pouze jednoho haplotypu (5/10) například všichni rodiče a děti. Umožňuje to podávání chemoterapie

pár dní po transplantaci, která zničí všechny buňky, které tělo nepřijme. Využívá se toho hlavně v případech časové tísně, kdy není čas hledat dárce v registrech. [3]

KIR geny se stejně jako HLA dědí celý blok. Jelikož HLA se nachází na 6 chromozomu a KIR na 19, tak shodní dárce v HLA znacích se jen menšinově shodují v KIR genech. V případě příbuzného dárce shodujícího se v HLA znacích je pouze 25% shodných také v KIR. [9]

TODO možná informaci o B-content score za tohle, proč u více shodných se řeší KIR, že se vybírají B haplotypový dárce a že v poslední době se řeší, zda kromě haplotypu nemají vliv konkrétní alelické varianty KIR genů (to je to, proč vy to žesíte v diplomce)

K zjištění konkrétních alelických variant se pro tzv. typizaci využívají sekvenační metody, typicky s polymerázovou řetězovou reakcí.

2.6 Bordel haplotypy

TODO budu to tam dopisovat? KIR2DL5 (Where two or more genes have very similar structures and have very similar sequences, they may be given the same number but distinguished by a final letter: for example, the KIR2DL5A and KIR2DL5B genes. The similarity of these two genes suggests they are related by a recent gene duplication event.),

3 Sekvenační metody získávání DNA dat

Po pojmem sekvence DNA se skrývá posloupnost písmen představujících primární strukturu reálné nebo hypotetické molekuly či vlákna DNA, které nese nějakou informaci. Jednotlivá písmena jsou označována jako nukleotidy nebo nukleové báze. Nukleové báze mohou být A - adenin, C - cytosin, G - guanin a T - thymin. [2]

Příkladem může být následující úsek sekvence na základě obrázku 2.1

ACGTCA (3.0.1)

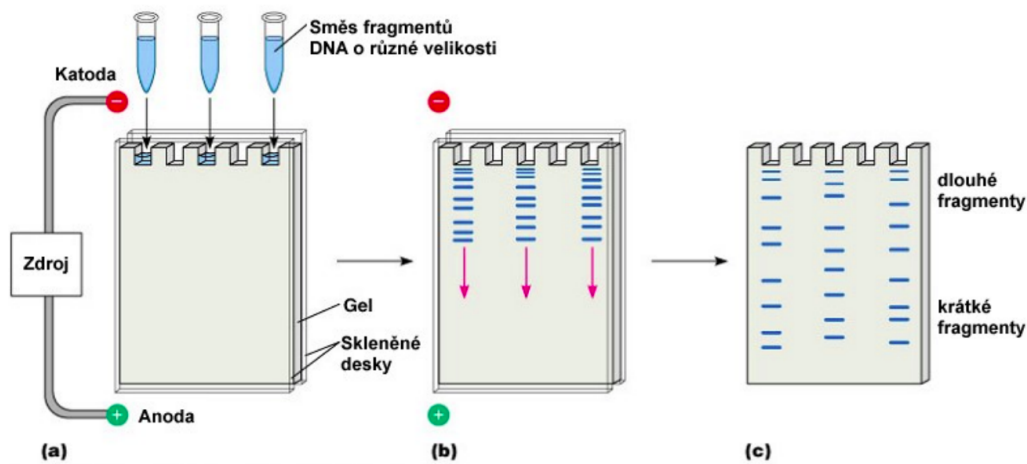
Sekvenování DNA, někdy pouze sekvenování, jsou biochemické metody, kterými se zjišťuje pořadí nukleotidů (A, C, G, T) v sekvenci DNA. Díky tomu je možné zjistit typizaci konkrétního člověka. Sekvenační metody se liší zejména délkou řetězce, kterou dokáží zpracovat, cenou a rychlostí sekvenace. Pro porovnání sekvenování celého genomu Sangerovo metodou by stálo několik milion dolarů a trvalo zhruba 10 let. Při použití dnešních metod by cena byla zhruba tisíc dolarů. Většina sekvenačních metod využívá vlastnosti přitahování báze do páru pouze jednou konkrétní bází. To znamená že se adenin vždy páruje s thyminem a cytosin se vždy páruje s guaninem. Z těchto párů vzniká již známá dvojité šroubovice DNA. Při sekvenování je možná se často setkat, že se sekvenuje jen konkrétní kus DNA, který je zrovna potřeba. Největším problémem u sekvenování je, že ready vzniklé ze sekvenátoru jsou jen kousky, které je třeba poskládat zpět. K tomu slouží zarovnávání. [15]

TODO možná tady ještě napsat něco o přípravě na sekvenování - je to dyžtak v té přednášce co nám říkala na FAV

3.1 Sanger sequencing

Sanger sekvenování využívá možnosti namnožení řetězce díky vzájemnému přitahování konkrétních bází. V prvním kroce replikace jsou nastříhané řetězce rozděleny na dvě vlákna. Lze si představit, že tyto dvě oddělená vlákna jsou dána do směsy, kde plavou jednotlivé nukleotydy spolu s upravenými

nukleotidy, které nesou specifickou fluorescenční barvu a za které není možné nic navázat. Následně za pomoci střídání teploty volně plující nukleotidy tvoří postupné páry s řetězcem, který chceme namnožit. Pokud se povede celý řetězec namnožit je odtržen a může se dále množit. Postupně ale bude docházet k navazování nukleotidů s fluorescenční barvou. Tím se vytvoří několik různě dlouhých sekvencí zakončených označeným nukleotidem. Podle jeho barvy je možné poznat o jaký nukleotid se jedná. Následně jsou za pomoci elektroforézy seřazeny v gelu podle délky. Elektroforéza rozděluje různě dlouhé sekvence na základě odlišnosti pohybu v elektrickém poli. Kratší doputují dále než delší. Pomocí sanger metody je možné sekvenovat řetězce dlouhé až 1000 bází.



Obrázek 3.1: Elektroforéza. [21]

3.2 NGS next-generation sekvenování

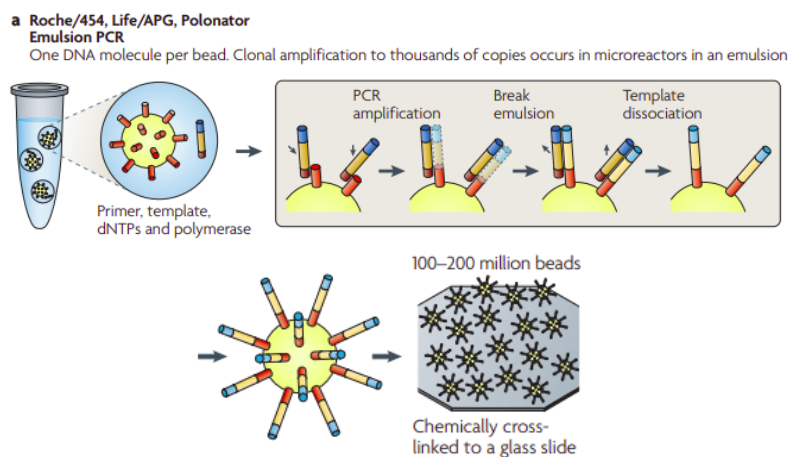
Next-generation sekvenování někdy označováno jako metody druhé generace jsou v porovnání se Sangerovo sekvenováním rychlejší a levnější, na druhou stranu ale dokáží zpracovávat jen řetězce dlouhé 100 až 500 bází, mají menší přesnost a častěji chybují. Jejich rychlost spočívá především ve schopnosti detekovat přidávání bází jednu po druhé a zároveň sekvenovat tisíce až miliony rozdílných molekul DNA najednou.

Všechny tyto metody si předpřipraví řetězce nastříháním na krátké části a připevním takzvaného adaptéru na jejich konec. Adaptér je krátká molekula DNA, která slouží k uchycení sekvenovaného úseku na pevný povrch. Řetězce DNA jsou namnoženy díky čemuž vzniknout klastry identických mo-

lekul koncentrovaných v jednom místě. Díky tomu je posílen signál, který by z pouhé jedné molekuly nebyl dostatečně silný. Tento signál je zachycen kamerou. Jeden z důvodů popularity NGS metod jsou i cenově dostupné stolní sekvenátory.

3.2.1 454 sekvenování a Ion Torrent

Pomocí 454 sekvenování je možné analyzovat více než milion molekul DNA najednou a délka každé jednotlivé sekvence se pohybuje okolo 700 až 1000 bází. V prvním kroku sekvenování je fragment DNA přichycena na malou "kuličku" na jejímž povrchu se postupně namnoží až kuličku zcela pokryjí identické fragmenty DNA. Následuje vložení kuličky i s DNA do jedné z milionů komůrek na destičce s reakční směsí. Postup znázorněn na obrázku 3.2. V určitém momentě je do této směsi přidán vždy jen jeden typ báze. Mezi jednotlivými fázemi přidávání určité báze jsou přebytečné nukleotidy z předešlého kroku odstraněny. To znamená že v reakční směsi je vždy jen jeden typ nukleotidů. Během vložení každé nové báze do rostoucího řetězce DNA je uvolněna molekula zvaná pyrofosfát, která spustí několik chemických reakcí. V poslední fázi enzym luciferáza vydá světelný záblesk, který je možné zachytit citlivou kamerou. Tento postup se nazývá pyrosekvenování. V případě, kdy je do řetězce přidáno několik stejných bází za sebou, například gen obsahuje podřetězec AAA, je vyzářeno, v našem případě, třikrát více světla než v případě jedné přiřazené báze. Kamera snímá celou destičku a na základě, která komůrka se rozsvítí pozná, kde proběhlo přidání báze. Intenzita světla pak určuje kolik bází bylo přidáno na jednu.



Obrázek 3.2: 454 sekvenování. [19]

Sekvenování Ion Torrent funguje na podobné principu sekvenování s roz-

dílem, že místo světla se měří změna pH v reakční směsi. Podle intenzity změny pH lze pak poznat kolik nukleotidů bylo přidáno do rostoucího řetězce.

Hlavní slabinou těchto dvou metod je značná chybovost při přidání mnoha stejných nukleotidů do řetězce za sebou. Například při přidání 10 A, nebude odpověď jednoznačná zda je to 10 A nebo 9.

3.2.2 Illumina

Při sekvenování pomocí Illumina jsou páry dvoušrobovice rozděleny na dva řetězce. Jednotlivé řetězce jsou následně přichyceny na malou destičku pomocí adaptéru. Každý řetězec se následně opakovaně množí až na destičce vznikne několik shluků. Přidání jedné molekuly ke druhé probíhá obdobně jako u Sanger sekvenování. Každý shluk tvoří jednu skupinu vzájemně identických řetězců. Mezi volné nukleotidy jsou opět zahrnuty nukleotidy označeny fluorescenční barvou za které nelze nic navázat. Oproti sangerovu sekvenování je ale tato blokáce vratná a po přečtení citlivou kamerou dojde k odstranění blokující části molekuly. Počítač si pak následně zpětně spočítá co to bylo za barvu (nukleotid). [5] [15]

3.2.3 SOLiD

SOLiD (Sequencing by Oligonucleotide Ligation and Detection) se spoléhá na enzym ligáza. Enzym je bílkovina, která určuje rychlost chemických reakcí. Enzym ligáza konkrétně umožňuje připojení jednořetězcových molekul k stávajícím řetězcům. K teplátu jsou přidávány takzvané sondy, což jsou kousky DNA. Sondy začínají všemi možnými dvojkombinacemi čtyř základních nukleotidů. V součtu je 16 sond. Na každé sondě je jedna ze čtyř fluorescenčních barev. V jednotlivých krocích jsou sondy připojeny k rostoucímu řetězci. Následně je přečtena fluorescenční barva, která je odstraněna a může se tak navázat další sonda. Z výsledného signálu lze pak odvodit sekvenci DNA.

3.3 Metody třetí generace

Velkým rozdílem oproti druhé generaci je že DNA templát není před sekvenováním namnožen a je čten pouze z jedné původní molekuly. Existuje například PacBio od Pacific Bioscience, který k detekci využívá fluorescenčně značené nukleotidy. Díky jeho vysoké citlivosti je možné v reálném čase zachytit

přidání i jediného nukleotidu do jediného řetězce DNA. Další zástupce je Oxford Nanopore jehož výhodou je jeho velikost. Oxford využívá odlišného tvaru bází. Obě metody jsou schopné přečíst přes 10 tisíc bází v rámci jedné analyzované molekuly DNA.

3.4 Read

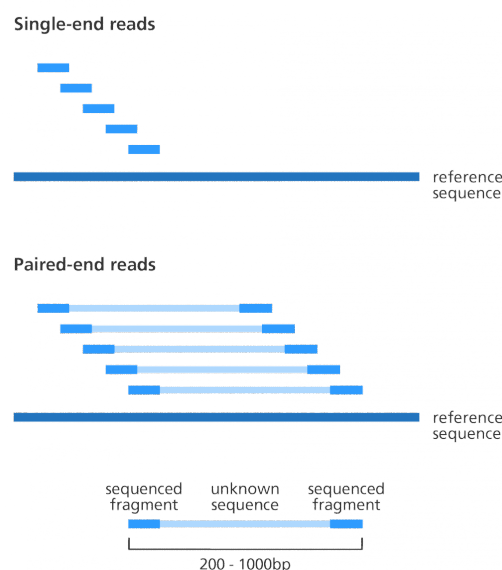
Read je sekvence bází odpovídající celému genomu či nějaké jeho části. Reads jsou typický výstup sekvenačních technik, kdy výstupem je sekvence nukleotidů o kterých nikdo neví co znamenají. Může to být gen, část genu nebo několik různých genů. Význam readu (o jaký gen se jedná) se zjišťuje zarovnáváním, kdy se daná sekvence porovnává vůči referenčnímu genu.

TODO je to z wiki, musím tady nutně udávat zdroj?

3.5 Single-end, paired-end a mate-pair

Single-end je sekvenování pouze jednoho konce molekuly. Nevýhoda tohoto způsobu se projeví především na krátkých readech, kde se zvýší problém jejich správného umístění. Oproti tomu v případě paired-end se sekvenuje z obou konci daného úseku. Vzniklé dva reads jsou označeny, v případě ART to naznačuje stejný název souboru spolu s 1 či 2 na jeho konci a zároveň je známá vzdálenost mezi oběma reads, která se pohybuje od 200 do 400 bp (base pair). Mate-pair je v podstatě paired-end s rozdílem, že je mezi reads větší vzdálenosti od 2 do 5 kb (kilobase) - takže přibližně 2000 - 5000 bp. [5]

TODO obrázek je z trochu blbího zdroje nejsem si jistá jestli ho můžu použít, ale mě přišel dobřej. <https://www.yourgenome.org/facts/how-do-you-put-a-genome-back-t>



Obrázek 3.3: Single-end a paired-end read.

3.6 Bordel

Sekvenování mRNA s použitím NGS technologií umožňuje měření genové exprese celého transkriptomu. Postup a provedení RNA-seq experimentu je znázorněn na obr. 14. Prvním úkolem je vyčistit zkoumaný vzorek o rRNA, tRNA a mitochondriální RNA, které u prokaryot i eukaryot tvoří přibližně 75 procent všech RNA molekul. Navzdory použití purifikačních metod, mezi které patří například poly(A)purifikace a DNS normalizace, sekvenční data mohou obsahovat menší množství těchto RNA molekul [59]. Ty mohou být odfiltrovány v následujících krocích bioinformatickými postupy. Zbýlá mRNA je poté nastříhána na menší části, a je z ní připravena knihovna krátkých fragmentů s navázanými adaptory. Ty jsou poté sekvenovány sekvenačním přístrojem a jako výsledek získáme tzv. ready. Samotné ready však nemají žádnou vypovídající hodnotu, a proto jsou dále bioinformaticky zpracovány. Namapováním na referenční sekvenci zjistíme jejich genomickou pozici, ze které byly odvozeny. Většina readů je namapována na exony, což jsou transkripčně aktivní jednotky, a pouze malé množství readů je namapováno na transposony. Ready které nejde namapovat v celku, jsou rozděleny na menší části a ty jsou namapovávány zvlášť. Rozdělené ready umožňují jednodušší identifikaci mezer mezi exony (angl. splice junctions) tohle je z té diplomky single-pair

4 Analyza dostupných bioinformatických nástrojů pro zpracování NGS dat

4.0.1 Vytvoření testovacího haplotypu

TODO navíc hned když se podívám na 3DL3: 00402 tak tam mám dvě možnosti TODO a co když tam nějaký není, tak prostě pokračuju tím dalším, nebo se tam něco dává jako mezera? Myslím když jste třeba napsala 2DL5B: TODO Co je 2DP1, ve FASTA jsem je nenašla tak jsem to skopčila z https://www.ebi.ac.uk/cgi-bin/ipd/kir/get_allele.cgi?2DP1*0020103 - p značí pseudogen A prej by měli být v KIR-gen.fasta ale já to tam za boha nemůžu najít jo tak tam je v tom KIR_gen fasta ale nahoře jsou i jiný co asi úplně nejsou pseudogeny tak to moc nechápu ..

to děláte správně, jen musíte vytvořit KIR haplotyp, který podšoupnete tomu ARTu, nejenom jeden konkrétní KIR, aby Vám vytvořil ready, tj. např. kombinaci (toplevel je jedna známá linie, sloučíte si ty KIRy za sebe):
3DL3: 00402, 00802 2DS2: 00101 2DL2: 00301 2DL3: 001 2DL5B: 2DS3:
2DP1: 00201 2DL1: 00302 3DP1: 007, 00901 2DL4: 00102, 00501 3DL1:
01502 3DS1: 01301 2DL5A: 001 2DS3: 2DS5: 00201 2DS1: 00201 2DS4: 001
3DL2: 0020105, 0070102

4.1 ART

ART (next-generation sequencing read simulator) je sada simulačních nástrojů, které generují syntetické ready, jako kdyby byli získány sekvenováním pomocí NGS. Nástroj ART dokáže simulovat single-end a paired-end ready ze sekvenátorů Illumina, 454 společnosti Roche a SOLid od společnosti Applied Biosystems. Ready, vytvořené nástrojem ART jsou používány pro testování a analýzy nástrojů zpracovávající právě NGS sekvence jako například zarovnávací (nástroj Bowtie). [13]

Podle [13] je dostupných několik simulačních nástrojů (Wgsim, MetaSim, SimSeq, FlowSim), které fungují dobře pro platformy pro které byly určeny, ale žádný z nich se nedokázal vypořádat se všemi hlavními platformami. Především v generování hlavních chyb, které jsou substituční a vložení či

smazání (INDEL - insert-deletion) na základě jednotlivých módů konkrétní platformy. ART obsahuje technologické profily chyb a navíc mu může být podsunut i uživatelský profil chyb. Obsažené profily délky readů a jejich chyb byly získány z datasetu skutečných sekvenovaných dat.

TODO možná někde zmínit co přesně znamená konkrétní chyba

TODO Proč? No protože přesně ví co tam dávají za data, protože mu podšoupnou ten referenční genom a tak pak můžou dobře sledovat co ten zarovnávač s tím dělá. A proč je to o tolik výhodnější než když by měli nějaký realnej dataset? Možná že si tam můžou ty chyby navolit tak jak se jim hodí? Jako bude v tom méně chyb, ale stejně.

TODO tohle někde dodat: ARTU je nutné podšoupnou referenční genom a ART vygeneruje ready způsobem napodobujícím sekvenční ready

Illumina je sekvenování založené na vratném umístění báze označené barvou do rostoucího řetězce jehož nejčastější chybou je substituce. Pravděpodobnost chyby substituce je určena na základě kvality skóre dané báze, které je závislé na pozici v rostoucím řetězci. Průměrné kvality skóre klesá v závislosti na zvyšování pozice báze. ART simuluje substituční chybu na základě tohoto skóre a empirického modelu získaného z trénovacích datasetů. INDEL chybu simuluje jen na základě empirického rozdělení z trénovacích dat. Pro paired-end simulaci, ART používá dvě rozdílné kvality skóre s distribucí a error rates pro první a druhý read.

454 je sekvenování při kterém se zachycuje vyzářené světlo na základě toho pokud se báze přidala do řetězce či nikoliv. Jeho dominantní chybou je tedy nesprávné určení počtů přidávaných bází. Pravděpodobnost chyby roste s frekvencí dlouhých úseků obsahujících stejnou bázi. Proto ART modeluje rozdělení chyb na základě délky úseku obsahující stejnou bázi spolu s Markovovy řetězcí.

SOLid je založené na označení čtyř barev pro 16 různých skupin bází. Pro paired-end read simulaci délky fragmentu je použito Gaussovské rozdělení. Rozdělení chyb je založeno na empirické znalosti získané z readů generovaných Applied Biosystémem. ART zároveň nabází nastavené chybovosti základě lineárního měřítka.

ART je implementován v jazyce C++ a je dostupný s licencí GPL verze 3

pro operační systémy Linux, MacOS a Windows. Je možné ho použít i jako C++ package. Pro jeho spuštění je nutné mít nainstalovaný compiler GNU g++ 4.0 nebo vyšší a knihovnu GNU gsl.

Data získána z FN Plzeň byla sekvenována nástrojem Illuminas proto i syntetické ready budou simulovat tento sekvenátor. Výstupy se čtou ve formátu FASTQ a zarovnání ve formátu ALN. může generovat zarovnávání také ve formátu SAM nebo UCS BED.

4.2 Bowtie

Bowtie je rychlý a paměťové efektivní nástroj pro zarovnávání krátkých sekvencí DNA na velké genomy. Bowtie je schopný zarovnat více než 25 milionů readů za hodinu (při běhu na jednom CPU) pro lidský genom s malým využitím paměti. Bowtie využívá FM indexaci s Burrows-Wheeler transformací a přidává k ní backtracking pro sledování nekonzistence. Novější verze Bowtie2 by měla být oproti Bowtie1 citlivější a rychlejší na delší ready než je 50 nukleotidů. Na lidský genom potřebuje 3.2 gigabajtů RAM. Nástroj bowtie je implementovaný v jazyce c++ s použitím knihovny SeqAn. [18] [17]

Bowtie je open source. V porovnání s nástroji Maq a SOAP je Bowtie rychlejší na lidském genomu. Citlivost má bowtie srovnatelnou s nástrojem SOAP a o něco menší než Maq. Ale je možnost pomocí příkazové řádky zvýšit citlivost na úkor rychlosti běhu programu. Bowtie vytváří indexy referenčních genů permanentní a lze je tak použít napříč běhy. Podporuje standardní vstupní formáty FASQ a FASTA. Výstupní zarovnání z bowtie je ve formátu SAM, což umožňuje návaznost s dalšími nástroji jako je třeba SAMtools

Zarovnávání bývá prvním krokem v mnoho genomických pipelinech. Často je to jejich nejpomalejší část, protože pro každý read musí zarovnávač vyřešit obtížný výpočetní problém. Určit pravděpodobné umístění v referenčním genomu. Mnoho zarovnávačů používá indexy k rychlému snižování kandidátů pro umístění zarovnávaného readu. FM index (Full-text minute-space) Přestože je tento vyhledávací prostor velký, mnoho jeho částí může být přeskočeno (odřezáno) bez ztráty citlivosti V praxi prořezávací strategie jako je dvojí indexování a obousměrné BWT usnadňuje v In practice, pruning strategies such as double indexing and bidirectional Burrows-Wheeler transform (BWT) facilitate very efficient untapped alignment of short reads.

Zarovnávání pomocí indexů může být i neefektivní v případě, že by bylo povoleno aby alignmenty obsahovali mezery. Mezery mohou být způsobeny chybou sekvenování nebo skutečným vložením a smazáním. Bezmezerován zarovnávaře, jako je bowtie, neodkážou zarovnat ready překlenující mezery a mezery zvyšují velikost vyhledávacího prostoru a zpomalují zarovnávání.

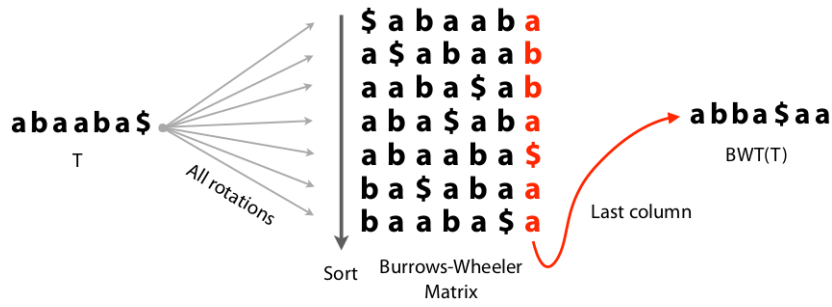
Pro každý read

1. extrahování seed z readů a jeho zpětné doplňky - nemyslí se tohle to paired atd
2. extrahované podřetězce jsou zarovnány na referenci v bezmezerové modelu za pomoci full-text minute index nebo neodděleně?
3. seed alignmenty jsou priorizovány a jejich pozice na referenčním genomu jsou spočítány z indexu
4. seedy jsou rozšířeny do úplného zarovnání pro zvýšení výkonu je použito SIMD -accelerated dynamic programming.

4.2.1 Burrows-Wheeler transformace

Burrows-Wheelerova transformace (BWT) je reverzibilní permutace řetězců v textu. Původně byla používána pro kompresy dat. Indexace založená na BWT umožňuje efektivní vyhledávání ve velké textu s malou pamětovou náročností.

BW transformace řetězce T , $BWT(T)$, je zobrazena na obrázku 4.1. Znak $\$$ je připojen na konec řetězce a zároveň musí platit, že se tento znak se v řetězci nevyskytuje. Burrows-Wheeler matice řetězce T je konstruovaná jako všechny cyklické rotace řetězce T , které byli seřazeny podle abecedy, kde znak $\$$ se bere, že je na začátku abecedy. Výstup, $BWT(T)$ pak představuje poslední sloupec matice. Tento řetězec má stejnou délku jako původní řetězec T . [18]



Obrázek 4.1: Burrows-Wheeler transformace řetězce T. [16]

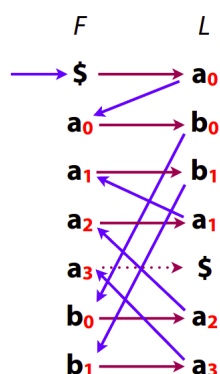
Burrows-Wheeler matice má vlastnost, která se nazývá last first mapping (LF). To znamená, že i -tý výskyt znaku X v prvním sloupci je i -tý výskyt znaku X v posledním sloupci. V případě přidání indexu do řetězce T je toto pravidlo pro znak a zobrazeno na obrázku 4.2. Obdobně to platí i pro ostatní znaky v řetězci.

$$T = a_0 b_0 a_1 a_2 b_1 a_3 \$ \quad (4.2.1)$$

F	L
\$	$a_0 b_0 a_1 a_2 b_1 a_3$
a_3	\$ $a_0 b_0 a_1 a_2 b_1$
a_1	$a_2 b_1 a_3$ \$ $a_0 b_0$
a_2	$b_1 a_3$ \$ $a_0 b_0$ a_1
a_0	$b_0 a_1 a_2 b_1 a_3$ \$
b_1	a_3 \$ $a_0 b_0 a_1$ a_2
b_0	$a_1 a_2 b_1 a_3$ \$ a_0

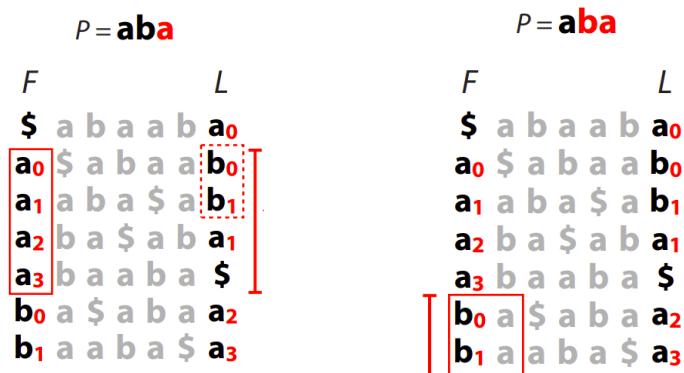
Obrázek 4.2: Burrows-Wheeler transformace last first mapping (LF). [16]

Zpětné získání řetězce je znázorněno na obrázku 4.3. L sloupec je řetězec který je výstupem BW transformace. F sloupec je snadné na základě L sloupce odvodit. Jelikož platí pravidlo, že počet jednotlivých znaků je stejný, stačí je pouze přemístit do F sloupce a seřadit podle abecedy. Dále s využitím LF je řetězec získán zpět. Jako první se vezme přidáný znak \$. Ve stejném řádku ve sloupci L se nachází a_0 . To znamená že řetězec začíná \$ a. Algoritmus pokračuje s a_0 v F sloupci. Ve stejném řádku v L sloupci je b_0 . b_0 je přidáno do řetězce a pokračuje až do doby než by byl opět znak \$.



Obrázek 4.3: Burrows-Wheeler transformace zpětné získání původního řetězce. [16]

Díky vztahu mezi F a L sloupcem je možné vyhledávat daný řetězec (zobrazeno na obrázku 4.4). Například vyhledávány řetězec bude $P = aba$. Při pohledu do F sloupce jsou nalezeny všechny sloupce začínající a , následně v L sloupci ve stejných řádcích jsou nalezeny dva výskyty b . Již je získán sufix ba , který existuje. Pokračuje se dále na řádky, které začínají právě nalezenými b . V sloupci L pro dané řádky jsou nalezena a . Řetězec $P = aba$ se v textu vyskytuje. FM index má více optimalizačních řešení, které dále nebudou rozebírány.



Obrázek 4.4: FM index - získání prefixu. [16]

Algoritmus, který by vyhledával přesné shody není v praxi použitelný, protože ready mohou obsahovat chyby vzniklé sekvenováním. Proto bowtie každé zarovnání zakládá na kvalitě znaku báze v daném readu. Bowtie postupně vytváří dlouhý sufix. Pokud se sufix nevyskytuje v textu pak se může algoritmus vrátit a v již vytvořeném sufixu nahradit bázi za jinou. Dále pokračuje obdobným způsobem. Pokud by měl algoritmus na výběr substituuovat za více bází vybere tu s nejnižší kvalitou znaku v readu. Protože bowtie

algoritmus je greedy je možné, že jeho nalezené řešení není to nejlepší. Pro nalezení toho nejlepšího řešení je třeba použít přepínač *--best*, jeho funkčnost je ale na úkor rychlosti, která může být 2x či 3x pomalejší. Zároveň je možné nastavit maximální počet nahrazených bází v readu.

4.2.2 bordel

It is particularly good at aligning reads of about 50 up to 100s or 1,000s of characters, and particularly good at aligning to relatively long (e.g. mammalian) genomes. Bowtie 2 supports gapped, local, and paired-end alignment modes.

Note that SOAP2 and Bowtie do not permit gapped alignment of unpaired reads.

Bowtie 2 by mělo být vhodnější pro delší ready než Bowtie1. We extracted a random subset of 1 million reads from each and aligned them with BWA-SW and Bowtie 2. We did not align with Bowtie, BWA or SOAP2 because those tools are designed for shorter reads. Bowtie už je překonanej nejenom Bowtie2 ale i BWA. Bowtie2 je podle studie znatelně lepší než Bowtie, SOAP2. tyhle výsledky jsou na syntetických readech

šla jsem přes docker docker image ls - zobrazí všechny image pak docker run a ID image sudo docker run -i -t 3c2b9a287f82 /bin/bash sudo docker ps -a

Tak jsem nakonec žádnéj docker nepotřebovala a stáhla jsem to tady po kliknutí na bowtie binary release.

na strance 25.4 je řečeno o hledání tch nejlepších zarovnání a je tam možnost *-best* ale že je dvakrát nebo třikrát pomalejší než normální mod.. a jde o to že najde první přijatelný a to označí kdežto při tom *best* prohledá co nejvíc a hledá to nejlepší i mezi těma přijatelnýma a to je pomalý.

takže zarovnání by mohlo být teoreticky namapování na referenční gen???

4.2.3 bordel

tak jsem to stáhla dala do složky a musela jsem teda nastavit proměnou prostředí export BT2_HOME=/home/kate/Dokumenty/FAV/Diplomka/existujicisw/bowtie2-2.4.1-linux-x86_64/ pak jsem pustila tohle: \$BT2_HOME/bowtie2-build \$BT2_HOME/example/reference/lambda_virus.falambda_virus a nakonec se mi vytvořili nějaký nové soubory lambda virus 1 atd.. v tom bowtie 2 adresáři

dělala jsem to podle tohohle webovky
z bowtie pak teda leze asi SAM formát

SAM

1. název readu který je zarovnáván

2. Sum of all applicable flags. Flags relevant to Bowtie are: součet všech aplikovaných (příslušných flags). Flagy relevantní k bowtie jsou: 1 - read je jeden z páru 2 - zarovnání je one z paired proper (The alignment is one end of a proper paired-end alignment) 4 - read má reported alignments 8 - read je jeden z páru a má reportovaný zarovnání 16 - zarování je obrácená reference vlákna 32 - The other mate in the paired-end alignment is aligned to the reverse reference strand 64 - read je mate 1 in a pair 128 - read je mate 2 in a pair

Thus, an unpaired read that aligns to the reverse reference strand will have flag 16. A paired-end read that aligns and is the first mate in the pair will have flag 83 ($= 64 + 16 + 2 + 1$).

3. jméno referencce ze které zarování patří 4. 1-based offset into the forward reference strand where leftmost character of the alignment occurs 1-based odszaneí v následující referenci 5. kvalita mapování 6. CIGAR reprezentace zarovnání 7. název reference kde je zarovnán kamarád 8. 1-based zarování ofsetu k následující referenci 9. Odvozená délka fragmentu. Velikost v závorku je že se mate nachází předtím. 0 že jsem nezarovnali mate 10. read sekvence 11. ASCII encoded read kvalita, stejné jako u FASTQ 12. optional pole

5 Bordel

5.1 ART

5.1.1 pokus to nejak spustit

Takze kdyz otebru hlavni readme tak mi to riká že tam jsou read me pro jednotlivy verze sekvenatoru ..

pak se to musí skompilovat

`./configure --prefix=$HOME make make install`

teď mě zajímá ta ilumina tak podle readme ilumina tak můžu vlést do složky examples a tam pustit skript `run_test_examples_illumina.sh`, tak tam jsou 4 příklady použití a pokud asi všechno dobře proběhne tak se mi zobrazí pár nových souborů ve složce examples..

FASTQ - *.fq data file s ready. pro paired-read simulator *1.fq obsahuje data pro první ready a *2.fq druhý ready

tohle nějak funguje MSv3 tam musím dát abych to mohla dostat na délku readu 250 a p znací ze to je paired.. tak se má používat MSv1 *art_illumina-ssMSv3-sam-iamplicon_reference.fa-p-l250-f10-m300-s10-omoje_art_data* Tohle používej: *art_illumina-ssMSv1-sam-iamplicon_reference.fa-p-l250-f100-m300-s10-omoje_art_data*

5.1.2 FASTQ

Sekvenační přístroje produkují data ve formátu FASTQ takže i ART musí logicky generovat tenhle formát. Pokud jsou ready v páru tak je na konci .1 a druhý read z páru tam má .2 to jsem u těch svých přímo nenašla

ale máš teda tři druhy single end, paired-end a matepair.

FASTQ obsahuje obě základy sekvence ?? both sequence bases a kvality skóre je to v následujícím formátu @read_id sequence read + base quality scores je kódovány by ascii code of a single character, kde je kvalita rovná score to ascii code character minus 33. chápu proč tam je to -33 protože když se podíváš do ascii tabulky tak je tam od 33 první normální znak jinak jsou tam divný .. takže třeba otazník je v ascii na 63 takže -33 takže má ohodnocení kvality 30 jen by mě teda zajímalo v jakém sme intervalu? - je 45 v ascii a nevím jestli to je teda od 0 do 100? a teda nejvyšší číslo znamená nejkvalitnější a nejmenší mín kvalitní? Podle té diplomky to tak je že čím vyšší číslo tím kvalitnější a většinou je to od 0 do 40 jen zřídka to překročí

hodnotu 60, když je tam 10 tak to znamená že jedna báze z deset je špatně.. když je tam 30 tak to znamená že jedna z 1000 je špatně. já tam mám třeba F a to je 70.

example: @refid-4028550-1 caacgccactcagcaatgatcggtttattcacgat... +

ALN - zarovnání readů zase *1.aln pro první a *2.aln pro druhý soubor je rozdělen na hlavičku a body část obsahuje hlavičku a v té hlavičce je jakým příkazem byl soubor vygenerován a reference na sequence id a jejich délku @CM tag pro příkaz a @SQ pro reference sequence Hlavička vždycky začíná s

HEADER EXAMPLE

v body jsou všechny zarovnání

aln_start_pos označuje počáteční pozici v referenci sekvenční, je vždy relativní vzhledem k vláknu referenční sekvenční To znamená že aln_start_pos plus (10) vlákno je odlišný od aln_start_pos minus (-) vlákna.. ??? WHAT???

ref_seq_aligned je zarovnaná oblast referenční sekvenční, která může být plus vlákno nebo mínus vlákno referenční sekvenční ref_seq_aligned je zarovnaný read, který je vždy ve stejné orientaci jako stejný read v odpovídajícím fastq suboru.

aln_start_pos is the alignment start position of reference sequence. aln_start_pos is always relative to the strand of reference sequence. That is, aln_start_pos 10 in the plus (+) strand is different from aln_start_pos 10 in the minus (-) stand.

ref_seq_aligned is the aligned region of reference sequence, which can be from plus strand or minus strand of the reference sequence. read_seq_aligned is the aligned sequence read, which always in the same orientation of the same read in the corresponding fastq file.

SAM je standardní formát pro NG sekvenční ready zarování BED o tom tam nic není jen NOTE: both ALN and BED format files use 0-based coordinate system while SAM format uses 1-based coordinate system.

pak jsou tady 4 doporučené použití *art_illumina[options]-ss < sequencing_system >*
-sam -i < seq_ref_file > -l < read_length > -f < fold_coverage >
-o < outfile_prefix > art_illumina[options] -ss < sequencing_system >
-sam -i < seq_ref_file > -l < read_length > -c < num_reads_per_sequence >
-o < outfile_prefix > art_illumina[options] -ss < sequencing_system >
-sam -i < seq_ref_file > -l < read_length > -f < fold_coverage >
-m < mean_frag_size > -s < std_frag_size > -o < outfile_prefix >
art_illumina[options] -ss < sequencing_system > -sam -i < seq_ref_file >
-l < read_length > -c < num_reads_per_sequence > -m < mean_frag_size >
-s < std_frag_size > -o < outfile_prefix >

pak tam máš parametry

a jak dlouhý chceme simulovat ready?

5.1.3 bordel

ART is freely available to public. The binary packages of ART are available for three major operating systems: Linux, Macintosh, and Windows. ART is also available as Platform-independent C++ source packages. Each package includes programs, documents and usage examples.

ART simuluje ready napodobobáním skutečných procesů sekvenování s empirickým chybovým modelem nebo quality profiles summarized from large recalibrated sequencing data ART může také simlovat čtené pomocí uživatelského vlastního read error modelu nebo quality profiles

TODO - tohle úplně nechápu ART podporuje simulaci jedno párových, dvou párových tří hlavních komerčních sekvenčních platfoem Výstupy se čtou ve formátu FASQ a zarování ve formátu ALN. ART může také generovat zarovnávání ve formátu SAM nebo UCSC BED ART lze použít společně se simulátory variant genomů VarSim

to je odtud 454 sekvenování je pyrosekvenování, které cycklicky testuje přítomnost každého ze čtyř nukleotidů DNA (T, A, C, G)

SOLid ke kódování 16 různých dinukleotidů používá čtyři fluoresenční barevná barviva, každé barvivo kóduje čtyři dinukleotidy

tak jsem stáhla normálně nejnovější verzi z niehs.nih.gov a podle instrukcí co byli v souboru INSTALL dala

musí se brát v potaz že z toho generátoru nikdy nebudou data taková jako reálná.. realná budou horší

SAM Sequence Alignment Map format), respektive jeho binárně komprimovaná verze BAM (z angl. Binary Alignment Map format).

5.2 IGV

nakonec jsem to pustila přes IGV ale stejně se tam museli ty indexi dodělat a musím být ve složce `/Dokumenty/FAV/Diplomka/existujicisw/IGV/IGV_Linux2.8.0`

když to otevřeš tak možná občas vypadá že tam nic není tak musíš vybrat konkrétní úsek nahore z toho rolóvátka Možná by se pak dalo udělat to že ty víš že tam může být maximálně dva z toho jednoho KIR souboru.. jako že může mít maximálně dvě alely z jednoho souboru No akorát co mi to udělá když bude mít dvě stejný?

6 Implementace

6.1 Popis problému

máme krátkou délku že read který dostáváme jsou 250 bp dlouhé a jeden gen může být dlouhý 14738 bp s tím že jednotlivé ready se nám tedy mohou překrývat můžou tam být chyby

teoreticky mám maximálně dvě možné alely s jednoho souboru, ale nemusím mít ani jednu

pak by se tam dala přidat heuristika že bych brala známe haplotypy

Možná pak ještě pracovat s pravděpodobností výskytu daného genu

možná by se pak ještě dalo kolik readů tam bylo zarovnaných- ale to je blbost protože tam mám ready z několika genů ne jen z toho jednoho

6.2 Návrh systému

6.3 Referenční geny

Referenční geny byly převzaty z IPD-KIR [24] konkrétně soubory ve formátu *fasta* uloženy ve stejnojmenné složce. Jednotlivé soubory jsou pojmenovány genem který obsahují např. *KIR2DL1_gen.fasta*. Každý soubor představuje všechny dostupné alely konkrétního genu.

Kromě souborů **_gen.fasta* obsahuje složka *fasta* také soubory **_prot.fast* a **_nuc.fasta*. Soubor **_prot.fast* obsahuje sekvenci proteinů, *nuc* obsahuje nucleotidy a *gen* obsahuje genomickou DNA sekvenci

TODO já vlastně nechápu jaký je rozdíl mezi *nuc* a *gen*?

All files in this folder are provided in the FASTA sequence format. Please note the FASTA format contains no alignment information.

Files designated “X_prot.fasta”, where X is a locus or gene, contain protein sequences. Please note that alleles that contain non-coding variations may be identical at the protein level.

Files designated “X_nuc.fasta”, where X is a locus or gene, contain the nucleotide coding sequences (CDS). Please note that alleles that contain non-coding variations may be identical at the CDS level.

Files designated “X_gen.fasta”, where X is a locus or gene, contain genomic DNA sequences. Please note for alleles that do not possess genomic sequences, there will be no entry in the file.

6.4 Použité programové prostředky

6.4.1 Python

Program byl navržen a implementován na operačním systému Linux za použití především programovacího jazyku Python. Pro spuštění programu je nutné mít nainstalovaný Python ve verzi 3.

Biopython

Biopython je sdružení vývojářů, kteří vytváří volně dostupné python nástroje vhodné pro výpočty v molekulární biologii. Biopython se snaží zjednodušit použití pythonu pro výzkum bioinformatiky. Mimo jiné umí pracovat s formáty souborů, které se využívají v bioinformatice jako je například BLAST nebo Fasta

Instalaci jde provést `pip install biopython`

tak to vypadá že i ten biopython umí aligned a že to dělá přes to Burrows wheeler aligner

7 Vyhodnocení výsledků a jejich srovnání

8 Závěr

9 Seznam zkratek

HLA KIR NK sekvence DNA cytosin, adenonin atd DNA (Deoxyribonukleovou kyselinu) co ty formáty souboru?

10 Výkladový slovník pojmů

WHO český národní registr možná zmínit národní registr genotyp fenotyp
tyhle kraviny Genotyp pro danou chromozomální oblast se pak u většiny lidí
skládá ze dvou haplotypů). genom kompletní sekvence daného organismu

Literatura

- [1] *Chromosome* [online]. [cit. 2020/12/3]. Dostupné z:
<https://www.genome.gov/genetics-glossary/Chromosome>.
- [2] *DNA sequencing Fact Sheet* [online]. [cit. 2019/03/1]. Dostupné z:
<https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Fact-Sheet>.
- [3] *S transplantací kostní dřeně stále častěji pomáhají příbuzní* [online].
Dostupné z: <https://ct24.ceskatelevize.cz/domaci/2527141-s-transplantaci-kostni-drene-stale-casteji-pomahaji-pribuzni>.
- [4] *Basic genetics* [online]. [cit. 2020/12/3]. Dostupné z:
<https://kintalk.org/genetics-101/>.
- [5] *Illumina* [online]. [cit. 2019/03/1]. Dostupné z:
<https://www.illumina.com/>.
- [6] *KIR genotypes* [online]. Dostupné z:
<http://www.allelefrequencys.net/kir6001a.asp>.
- [7] BARANWAL, A. – MEHRA, N. Major Histocompatibility Complex Class I Chain-Related A (MICA) Molecules: Relevance in Solid Organ Transplantation. *Frontiers in Immunology*. 02 2017, 8. doi: 10.3389/fimmu.2017.00182.
- [8] BERNAREGGI, D. – POUYANFARD, S. – KAUFMAN, D. S. Development of innate immune cells from human pluripotent stem cells. 2019. Dostupné z:
<https://www.sciencedirect.com/science/article/pii/S0301472X19300037?via%3Dihub>.
- [9] COOLEY, S. – WISDORF, D. J. – GUETHLEIN, L. A. Donor selection for natural killer cell receptor genes leads to superior survival after unrelated transplantation for acute myelogenous leukemia. 2010. Dostupné z:
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2953880/#>.
- [10] FRYČOVÁ, M. Lze u pacientů s AML indikovaných k nepříbuzenské transplantaci provádět v klinické praxi výběr nepříbuzných dárců na základě KIR genotypů, 2016.
- [11] HERNYCHOVÁ, L. Receptory NK buněk. 2012.

- [12] HSU, K. C. et al. *The killer cell immunoglobulin-like receptor (KIR) genomic region: gene-order, haplotypes and allelic polymorphism* [online]. 2002. Dostupné z: <https://onlinelibrary.wiley.com/doi/full/10.1034/j.1600-065X.2002.19004.x>.
- [13] HUANG, W. et al. ART: a next-generation sequencing read simulator. 2012. Dostupné z: <https://academic.oup.com/bioinformatics/article/28/4/593/213322>.
- [14] J, R. et al. *Nomenclature* [online]. Nucleic Acids Research, 2015. [cit. 2019/10/1]. 43:D423-431. Dostupné z: <http://hla.alleles.org/misc/citing.html>.
- [15] KOLÍSKO, M. Moderní metody sekvenování DNA. 2017. Dostupné z: <https://ziva.avcr.cz/files/ziva/pdf/moderni-metody-sekvenovani-dna.pdf>.
- [16] LANGMEAD, B. [online]. [cit. 2019/03/1]. Dostupné z: <http://www.langmead-lab.org/>.
- [17] LANGMEAD, B. – SALZBERG, S. L. Fast gapped-read alignment with Bowtie 2. 2012. Dostupné z: <https://www.nature.com/articles/nmeth.1923>.
- [18] LANGMEAD, B. et al. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. 2009. Dostupné z: <https://genomebiology.biomedcentral.com/articles/10.1186/gb-2009-10-3-r25>.
- [19] MERZKER, M. L. Sequencing technologies-the next generation. 2013. doi: 10.1038/nrg2626.
- [20] MUDR. PAVEL JINDRA, P. D. *Imunopatologické a imunogenetické aspekty transplantací krvetvorných buněk a solidních orgánů*. PhD thesis, Universita Karlova v Praze, 2011.
- [21] PAPOUŠEK, I. Elektroforéza nukleových kyselin. 2017. Dostupné z: https://fvhe.vfu.cz/files/mbhp_2018_02.pdf.
- [22] PENKA, M. – KOLEKTIV, E. T. *Hematologie a transfuzní lékařství II*. 2012. ISBN 978-80-247-3460-6.
- [23] ROBINSON, J. et al. IPD—the Immuno Polymorphism Database. 2013. Dostupné z: <https://www.ebi.ac.uk/ipd/index.html>.
- [24] ROBINSON, J. et al. The IPD and IMGT/HLA Database:allele variant databases. 2015. Dostupné z: <https://www.ebi.ac.uk/ipd/index.html>.

- [25] S.KANNANA, G. – ARIANEXYSAQUINO-LOPEZ – A.LEED, D. Natural killer cells in malignant hematology: A primer for the non-immunologist. 2017. Dostupné z: <https://www.sciencedirect.com/science/article/pii/S0268960X16300704>.
- [26] SMITH, D. T. *Encyklopedie lidského těla*. 2005. ISBN 80-7321-156-4.
- [27] THIELENS, A. – VIVIER, E. – ROMAGNÉ, F. NK cell MHC class I specific receptors (KIR): from biology to clinical intervention. *Current opinion in immunology*. 2012, 24 2, s. 239–45.

A Uživatelská dokumentace

A.1 Nastavení ART a jeho spuštění